

DIALOG(R)File 351:Derwent WPI
(c) 2003 Thomson Derwent. All rts. reserv.

011404945 **Image available**

WPI Acc No: 1997-382852/ 199735

XRPX Acc No: N97-318669

Space detection method for documents containing both English and Japanese characters - by evaluating length of extracted object range to determine if its length is constant or proportional

Patent Assignee: RICOH KK (RICO)

Number of Countries: 001 Number of Patents: 001

Patent Family:

Patent No	Kind	Date	Applicat No	Kind	Date	Week
JP 9167206	A	19970624	JP 95328737	A	19951218	199735 B

Priority Applications (No Type Date): JP 95328737 A 19951218

Patent Details:

Patent No	Kind	Lan Pg	Main IPC	Filing Notes
JP 9167206	A	13		

Abstract (Basic): JP 9167206 A

The method involves extracting an object range consisting of several English character rows. The length of the extracted object range is evaluated when performing a length format evaluation process (301) to determine if it is a proportional length or a constant length.

When the length of the object range is constant, a constant length character string detection process (303) detects a space for every object range. When the length of the object range is proportional, a proportional length character string detection process (302) detects a space for every object range.

ADVANTAGE - Provides accurate space detection by using length of extracted character string.

Dwg.1/9

Title Terms: SPACE; DETECT; METHOD; DOCUMENT; CONTAIN; ENGLISH; JAPAN; CHARACTER; EVALUATE; LENGTH; EXTRACT; OBJECT; RANGE; DETERMINE; LENGTH; CONSTANT; PROPORTION

Derwent Class: T01; T04

International Patent Class (Main): G06K-009/36

International Patent Class (Additional): G06K-009/62

File Segment: EPI

Manual Codes (EPI/S-X): T01-J10B2; T04-D03; T04-D04

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-167206

(43) 公開日 平成9年(1997)6月24日

(51) Int.Cl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 K 9/36			G 0 6 K 9/36	
9/62	6 1 0	9061-5H	9/62	6 1 0 B

審査請求 未請求 請求項の数 7 O L (全 13 頁)

(21) 出願番号 特願平7-328737

(22) 出願日 平成7年(1995)12月18日

(71) 出願人 000006747

株式会社リコー

東京都大田区中馬込1丁目3番6号

(72) 発明者 嶺脇 隆邦

東京都大田区中馬込1丁目3番6号 株式会社リコー内

(74) 代理人 弁理士 鈴木 誠 (外1名)

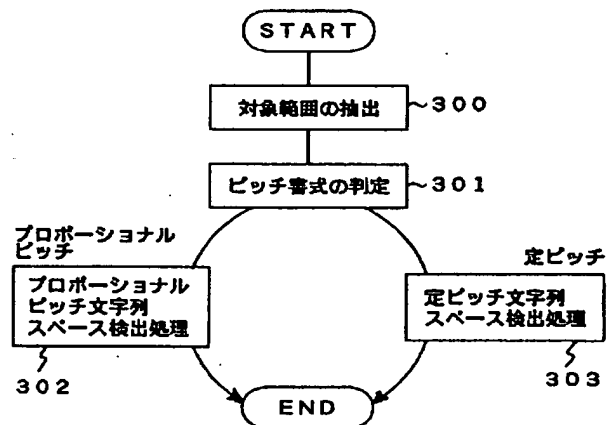
(54) 【発明の名称】 日英混在文書のスペース検出方法、ピッチ書式判定方法、定ピッチ英数文字列のスペース検出方法、及びプロポーショナルピッチ英数文字列のスペース検出方法

(57) 【要約】

スペース検出処理フロー

【課題】 日本文字列と混在した英数文字列中のスペースを高精度に検出する。

【解決手段】 文字認識装置において、認識結果から英数文字列を抽出し(300)、英数文字列毎にピッチ書式を判定する(301)。各英数文字列に対し、そのピッチ書式に応じプロポーショナルピッチ用スペース検出処理(302)又は定ピッチ用スペース検出処理(303)を適用し、高精度にスペースを検出する。



【特許請求の範囲】

【請求項1】 文書の画像より文字を切り出し認識する文字認識装置において、英日混在の文書の文字認識結果に基づいて英数文字列を対象範囲として抽出する対象範囲抽出処理と、該対象範囲抽出処理により抽出された各対象範囲毎にピッチ書式が定ピッチかプロポーショナルピッチかを判定するピッチ書式判定処理と、該ピッチ書式判定処理により定ピッチと判定された各対象範囲毎にスペースを検出する定ピッチ文字列スペース検出処理と、該ピッチ書式判定処理によりプロポーショナルピッチと判定された各対象範囲毎にスペースを検出するプロポーショナルピッチ文字列スペース検出処理とを有することを特徴とする日英混在文書のスペース検出方法。

【請求項2】 文書の画像より文字を切り出し認識する文字認識装置において、英日混在の文書の文字認識結果に基づいて英数文字列を対象範囲として抽出する対象範囲抽出処理と、該対象範囲抽出処理により抽出された各対象範囲毎にピッチ書式が定ピッチかプロポーショナルピッチかを判定するピッチ書式判定処理と、該ピッチ書式判定処理により定ピッチと判定された各対象範囲毎にスペースを検出する定ピッチ文字列スペース検出処理と、該ピッチ書式判定処理によりプロポーショナルピッチと判定された各対象範囲毎にスペースを検出するプロポーショナルピッチ文字列スペース検出処理とを有し、該ピッチ書式判定処理において、各対象範囲に対し、各対象範囲毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が所定の閾値より小さい文字矩形間隔のほうが、文字矩形間隔比が該閾値以上の文字矩形間隔より多数であるならばプロポーショナルピッチと判定し、そうでなければ定ピッチと判定することを特徴とする日英混在文書のスペース検出方法。

【請求項3】 文書の画像より文字を切り出し認識する文字認識装置において、英日混在の文書の文字認識結果に基づいて英数文字列を対象範囲として抽出する対象範囲抽出処理と、該対象範囲抽出処理により抽出された各対象範囲毎にピッチ書式が定ピッチかプロポーショナルピッチかを判定するピッチ書式判定処理と、該ピッチ書式判定処理により定ピッチと判定された各対象範囲毎にスペースを検出する定ピッチ文字列スペース検出処理と、該ピッチ書式判定処理によりプロポーショナルピッチと判定された各対象範囲毎にスペースを検出するプロポーショナルピッチ文字列スペース検出処理とを有し、該定ピッチ文字列スペース検出処理において、注目した文字間の前後の文字間の文字矩形ピッチのうちの小さい方の文字矩形ピッチを注目した文字間の基準ピッチとし、注目した文字間の文字矩形ピッチの基準ピッチとの比が所定の閾値より大きいときに、注目した文字間にスペースが存在すると判定することを特徴とする日英混在文書のスペース検出方法。

【請求項4】 文書の画像より文字を切り出し認識する文字認識装置において、英日混在の文書の文字認識結果に基づいて英数文字列を対象範囲として抽出する対象範囲抽出処理と、該対象範囲抽出処理により抽出された各対象範囲毎にピッチ書式が定ピッチかプロポーショナルピッチかを判定するピッチ書式判定処理と、該ピッチ書式判定処理により定ピッチと判定された各対象範囲毎にスペースを検出する定ピッチ文字列スペース検出処理と、該ピッチ書式判定処理によりプロポーショナルピッチと判定された各対象範囲毎にスペースを検出するプロポーショナルピッチ文字列スペース検出処理とを有し、該プロポーショナルピッチ文字列スペース検出処理において、各対象範囲毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が、所定の閾値より大きいときに、対応する文字間にスペースが存在すると判定することを特徴とする日英混在文書のスペース検出方法。

【請求項5】 文書の画像より文字を切り出し認識する文字認識装置において、認識された各英数文字列に対し、各英数文字列毎に切り出し情報に基づいて決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が所定の閾値より小さい文字矩形間隔のほうが、文字矩形間隔比が該閾値以上の文字矩形間隔より多数であるならばプロポーショナルピッチと判定し、そうでなければ定ピッチと判定することを特徴とするピッチ書式判定方法。

【請求項6】 文書の画像より文字を切り出し認識する文字認識装置において、認識された定ピッチの各英数文字列において、注目した文字間の前後の文字間の文字矩形ピッチのうちの小さい方の文字矩形ピッチを注目した文字間の基準ピッチとし、注目した文字間の文字矩形ピッチの基準ピッチとの比が所定の閾値より大きいときに、注目した文字間にスペースが存在すると判定することを特徴とする定ピッチ英数文字列のスペース検出方法。

【請求項7】 文書の画像より文字を切り出し認識する文字認識装置において、認識されたプロポーショナルピッチの各英数文字列において、各英数文字列毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が、所定の閾値より大きいときに、対応する文字間にスペースが存在すると判定することを特徴とするプロポーショナルピッチ英数文字列のスペース検出方法。

【発明の詳細な説明】**【0001】**

【発明の属する分野】 本発明は、文字認識装置におけるスペース検出処理に係り、特に、日英混在の文書に対するスペース検出処理に関する。

【0002】

【従来の技術】 文書を文字認識によってテキストデータ

に変換し、このテキストデータから文書を再現できるようにする場合等には、文字認識に際し、文書中のスペースも検出する必要がある。このような文書中のスペース検出に関する従来技術としては次に述べるようなものが知られている。

【0003】従来技術1：文字間の白画素数を計数し、計数値を標準文字ピッチで割ることによりスペースコード数を決定する（特願昭63-14282号）。

【0004】従来技術2：文書のピッチ書式（定ピッチ／プロポーショナルピッチ）を判定する。そして、判定したピッチ書式に応じた方法でスペースを検出する（特願平1-161176号）。

【0005】従来技術3：日英文字混在の日本語文章において、文字間空白幅と文字幅の情報を比較しスペースの種類と有無を判定するが、スペースの前後にある文字の種類を参照してスペースの有効／無効を判定する（特願平2-214136号）。

【0006】従来技術4：日英文字混在の日本語文章において、文字間空白幅と文字の標準サイズを比較し、スペースの種類と有無を判定するが、文字間空白に隣接する文字に応じて文字間空白幅を補正し、この補正後の値を標準文字サイズと比較する（特願平3-18476号）。

【0007】

【発明が解決しようとする課題】日本文字と英数字（英文字と数字）が混在する文書では、日本文字と英数字のフォントの違い、全角文字と半角文字という文字サイズの違い、定ピッチとプロポーショナルピッチというピッチ書式の違いが存在し、しかも、このような違いが文書中の不特定の部分に生じる。そして、スペースの幅も場所によって違いが生じる。このことは図9に示す印字サンプルをみれば容易に理解されよう。図9において、全て全角文字で定ピッチ印字されたサンプル1では、「R i c o h」の「h」と次の「P」との間にだけスペースがある。これに対して英数字がプロポーショナルピッチで印字されたサンプル2では、「新型」の「型」と次の「R」との間、及び「R i c o h」の「h」と次の「P」にスペースがあり、しかも、そのスペースはサンプル1のスペースより間隔が狭い。

【0008】このような日英混在文書に対しては、従来技術では高精度のスペース検出が困難な場合があった。また、行中で日本文字と英数字が混在するような場合にピッチ書式を的確に判定する方法は知られていなかった。従来技術2は、文書全体又は行全体が英文であると仮定し、英単語単位の処理となっているため、行中に部分的に現れる英単語中のスペースを精度よく検出できなかった。定ピッチ文字列中のスペース検出に関しては、定ピッチ文字列の文字間隔の変動が大きいため、全角サイズのスペースは検出できても、文字間隔を補正する従来技術4によっても半角スペースの検出精度が上がらな

かった。プロポーショナルピッチ文字列中のスペース検出に関しては、プロポーショナルピッチ文字列の文字間隔が狭いので、従来技術4によってもスペースを検出できなかった。

【0009】本発明の目的は、日本文字と英数字が行中に混在するような日英混在文書に対し高精度のスペース検出が可能な改良した方法、並びに、このスペース検出方法のために好適なピッチ書式の判定方法、定ピッチ文字列中のスペースの検出方法、及びプロポーショナルピッチ文字列中のスペースの検出方法を提供することにある。

【0010】

【課題を解決するための手段】請求項1記載の発明による日英混在文書のスペース検出方法は、文書の画像より文字を切り出し認識する文字認識装置において、英日混在の文書の文字認識結果に基づいて英数字文字列を対象範囲として抽出する対象範囲抽出処理と、該対象範囲抽出処理により抽出された各対象範囲毎にピッチ書式が定ピッチかプロポーショナルピッチかを判定するピッチ書式判定処理と、該ピッチ書式判定処理により定ピッチと判定された各対象範囲毎にスペースを検出する定ピッチ文字列スペース検出処理と、該ピッチ書式判定処理によりプロポーショナルピッチと判定された各対象範囲毎にスペースを検出するプロポーショナルピッチ文字列スペース検出処理とを有することを特徴とするものである。

【0011】請求項2記載の発明は、請求項1記載の発明の日英混在文書のスペース検出方法において、ピッチ書式判定処理で、各対象範囲に対し、各対象範囲毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が所定の閾値より小さい文字矩形間隔のほうが、文字矩形間隔比が該閾値以上の文字矩形間隔より多数であるならばプロポーショナルピッチと判定し、そうでなければ定ピッチと判定することを特徴とするものである。

【0012】請求項3記載の発明は、請求項1記載の発明の日英混在文書のスペース検出方法において、定ピッチ文字列スペース検出処理で、注目した文字間の前後の文字間の文字矩形ピッチのうちの小さい方の文字矩形ピッチを注目した文字間の基準ピッチとし、注目した文字間の文字矩形ピッチの基準ピッチとの比が所定の閾値より大きいときに、注目した文字間にスペースが存在すると判定することを特徴とするものである。

【0013】請求項4記載の発明は、請求項1記載の発明の日英混在文書のスペース検出方法において、プロポーショナルピッチ文字列スペース検出処理で、各対象範囲毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が、所定の閾値より大きいときに、対応する文字間にスペースが存在すると判定することを特徴とするものである。

【0014】請求項5記載の発明のピッチ書式判定方法

は、文書の画像より文字を切り出し認識する文字認識装置において、認識された各英数文字列に対し、各英数文字列毎に切り出し情報に基づいて決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が所定の閾値より小さい文字矩形間隔のほうが、文字矩形間隔比が該閾値以上の文字矩形間隔より多数であるならばプロポーショナルピッチと判定し、そうでなければ定ピッチと判定することを特徴とする。

【0015】請求項6記載の発明の定ピッチ英数文字列のスペース検出方法は、文書の画像より文字を切り出し認識する文字認識装置において、認識された定ピッチの各英数文字列において、注目した文字間の前後の文字間の文字矩形ピッチのうちの小さい方の文字矩形ピッチを注目した文字間の基準ピッチとし、注目した文字間の文字矩形ピッチの基準ピッチとの比が所定の閾値より大きいときに、注目した文字間にスペースが存在すると判定することを特徴とするものである。

【0016】請求項7記載の発明のプロポーショナルピッチ英数文字列のスペース検出方法は、文書の画像より文字を切り出し認識する文字認識装置において、認識されたプロポーショナルピッチの各英数文字列において、各英数文字列毎に切り出し情報に基づき決定される標準文字サイズに対する文字矩形間隔の比である文字矩形間隔比が、所定の閾値より大きいときに、対応する文字間にスペースが存在すると判定することを特徴とするものである。

【0017】

【発明の実施の形態】本発明の実施の形態を明らかにするため、図面を用いて本発明の一実施例を説明する。図1は本発明の一実施例の全体的処理フローを示し、図2は本発明の一実施例のための装置構成例を示す。図1中のスペース検出処理のフローを図3に示す。図3中のピッチ書式判定処理のフローを図4に、プロポーショナルピッチ文字列スペース検出処理のフローを図5に、定ピッチ文字列スペース検出処理のフローを図6に、それぞれ示す。図7及び図8はスペース検出処理の説明のための図である。

【0018】初めに図1及び図2を参照し、全体的処理フローを装置構成と関連付けて説明する。まず、画像入力部200により文書画像のデータを入力し、これをバス213を經由して画像メモリ208に格納する(ステップ100)。画像が入力されると、行・文字切り出し部202において、画像メモリ208内の文書画像中の文字行と文字を切り出し、切り出した文字行及び文字の領域の情報を切り出しメモリ209に格納する(ステップ110)。ここに格納される切り出し情報は、例えば、文字行の始点と終点の座標あるいは始点座標と幅、文字の外接矩形(文字矩形)の対角頂点の座標やサイズ情報である。これらの情報は、各文字矩形と、それが所属する文字行との対応関係が識別できるような形で作成

され格納されることは当然である。そして、文字認識部203において、切り出しメモリ209内の切り出し情報を参照し、画像メモリ208より文字画像を取り込み、それを文字辞書メモリ210内の文字辞書と比較することにより文字画像に対する文字コードを決定し、それを結果メモリ211に格納する(ステップ120)。なお、文字認識部203において、単語や文法等の言語知識を利用して、文字辞書との比較による認識結果に対する修正等の後処理を行ってもよい。

【0019】文字認識部203による認識処理が文書画像の全体あるいは一部について終了した段階で、スペース検出処理(ステップ130)が開始する。このスペース検出処理(ステップ130)は、対象範囲抽出部204、ピッチ書式判定部205、定ピッチ文字列スペース検出部206及びプロポーショナルピッチ文字列スペース検出部207により実行され、その際に切り出しメモリ209、結果メモリ211及びワークメモリ212が参照される。スペース検出結果は結果メモリ211に格納される。スペース検出処理が終了すると、結果メモリ211内のデータが外部に出力され(ステップ140)、処理全体が終了する。

【0020】スペース検出処理(ステップ130)の概略は図3に示すとおりである。まず、対象範囲抽出部204において、結果メモリ211内の文字認識結果(文字コード)を参照し、スペース検出処理の対象範囲としての英数文字列を抽出する(ステップ301)。抽出した対象範囲に関する情報はワークメモリ212に保存される。次に、ピッチ書式判定部205において、対象範囲の文字に関して切り出しメモリ209及び結果メモリ211の内容を参照し、対象範囲がプロポーショナルピッチであるか、定ピッチであるかを判定する(ステップ301)。判定結果はワークメモリ212に保存される。プロポーショナルピッチと判定された対象範囲に対しては、プロポーショナルピッチ文字列スペース検出部207において切り出しメモリ209及び結果メモリ211の内容を参照して、プロポーショナルピッチ用の方法によりスペース検出を行い、検出したスペースをスペースコードとして結果メモリ211に書き込む(ステップ302)。定ピッチと判定された対象範囲に対しては、定ピッチ文字列スペース検出部206において切り出しメモリ209及び結果メモリ211の内容を参照して、定ピッチ用の方法によりスペース検出を行い、検出したスペースをスペースコードとして結果メモリ211に書き込む(ステップ303)。こここ述べた処理の具体的な内容について、図4乃至図8を参照し説明する。

【0021】図3の(a)に示すような文字行を処理する場合を考える。この文字行には、スペースを含む部分が2カ所ある。一つはプロポーショナルピッチ(フォントはTimes-Roman)で印字された「Preter 5」の「Preter」と「5」の間である。もう一

つは、定ピッチ（フォントはCourier）で印字された「Imagio 77」の「Imagio」と「77」の間である。この文字行に対する文字認識（ステップ120）の結果はスペースを含まず、「新型コピー、Preter55とImagio77発売」となる。

【0022】対象範囲抽出処理（ステップ300）では、そのようなスペースを含まない文字認識結果を参照し、英数字（英字又は数字）が5文字以上連続している文字列を対象範囲として抽出する。したがって、図7の（a）に示す文字行では、図7の（b）に示すように「Preter55」と「Imagio77」の二つの文字列が対象範囲として抽出される。なお、抽出判定のための文字長は5文字に限定されるものではなく、処理する文書の種類に応じて適宜変更してよい。また、対象範囲の決定方法そのものも適宜変更し得るものである。

【0023】このようにして対象範囲として抽出された英数文字列について、ピッチ書式判定（ステップ301）を行う。図4はそのフローチャートである。

【0024】まず、対象範囲に関する標準文字サイズを決定する（ステップ400）。具体的には、対象範囲内における最大の文字矩形高さの1.25倍（文字サイズA）と、対象範囲が含まれる文字行内の最大の文字矩形高さ（文字サイズB）を求める。文字サイズA、Bの大きい方の値を標準文字サイズとする。ただし、対象範囲内に「j」の文字矩形がある場合、その文字矩形は標準文字サイズの決定には利用しない。図8は、このような標準文字サイズ決定の説明図である。

【0025】図7に示した一つ目の対象範囲では、「P」の文字矩形の高さが最大であるので、その高さを1.25倍した値と、文字行内で最大の「新」の文字矩形高さとを比較し、大きい方の値が標準文字サイズとされる。なお、他の方法によって標準文字サイズを決定してもよい。

【0026】以上のようにして対象範囲に対する標準文字サイズを決定すると、次に、この標準文字サイズを用いて、対象範囲の先頭より、隣接する文字矩形のペアを選び、その文字矩形間の空白部の幅、すなわち文字矩形間隔（図7の（c）参照）の標準文字サイズに対する比（文字矩形間隔比＝文字矩形間隔／標準文字サイズ）を計算する（ステップ405、410）。計算した文字矩形間隔比は、ワークメモリ212に保存される。文字矩形間隔は、切り出しメモリ209に保存されている二つの文字矩形の座標から算出される。ただし、切り出し処理の段階で文字矩形間隔を求めておいてもよく、その場合は、切り出しメモリ209より読み出した文字矩形間隔を文字矩形間隔比の計算に用いればよい。次に、このようにして計算した文字矩形間隔比を所定の閾値TH1と比較する（ステップ415）。この閾値TH1は例えば0.17である（この値に限らないが、一般的文書では、このあたりの値が適当であることが実験により確認

された。ただし、処理する文書に応じて調整するのが好ましい）。文字矩形間隔比が閾値TH1より小さいときにカウンタNS（初期値は0）を1だけインクリメントし（ステップ420）、文字矩形間隔比が閾値TH1以上のときにカウンタNL（初期値は0）を1だけインクリメントする（ステップ421）。次の文字矩形ペアが残っている場合は、ステップ405からステップ420の処理ループを繰り返す。

【0027】最後の文字矩形ペアまで処理が終了すると、ステップ405からステップ420又は421までの処理ループを抜けステップ425に進む。このステップ425では、カウンタNSの値とカウンタNLの値を比較する（ステップ425）。カウンタNSの値は対象範囲内における文字矩形間隔比が閾値TH1より小さい、つまり「標準文字サイズ×TH1」より狭い文字矩形間隔の個数であり、カウンタNLの値は広い文字矩形間隔の個数である。ゆえに、ステップ425では、狭い文字矩形間隔が過半数であるか否かを調べている。NS>NLつまり狭い文字矩形間隔が過半数と判定した場合には対象範囲の文字列のピッチ書式をプロポーショナルピッチに設定し、そのフラグ情報をワークメモリ212に対象範囲と対応付けて書き込み（ステップ430）、そうでない場合、つまり広い文字矩形間隔が半数以下である場合には、ピッチ書式を定ピッチに設定し、そのフラグ情報をワークメモリ212に対象範囲と対応付けて書き込む（ステップ435）。このフラグ情報にしたがって、次の処理として、プロポーショナルピッチ文字列スペース検出処理（ステップ302）又は定ピッチ文字列スペース検出処理（ステップ303）が選択される。

【0028】図7の（a）に示した例では、最初の対象範囲「Preter55」は図7（c）に示すようにプロポーショナルピッチと判定され、もう一つの対象範囲「Imagio77」は定ピッチと判定される。

【0029】プロポーショナルピッチと判定された対象範囲に対するスペース検出処理（ステップ302）の処理内容について説明する。図5は、そのフローチャートである。

【0030】ワークメモリ212に保存されている対象範囲の文字矩形間隔比を取り出し、それを所定の閾値TH2と比較する（ステップ500、505）。この閾値TH2は例えば0.29である（この値に限らないが、一般的文書では、このあたりの値が適当であることが実験により確認された。ただし、処理する文書に応じて調整するのが好ましい）。文字矩形間隔比がTH2以下であれば、つまり文字矩形間隔が「標準文字サイズ×TH2」より狭いときには、対応した文字矩形ペアの間にスペースは無いと判断される。文字矩形間隔比がTH2より大きいときには、つまり文字矩形間隔が「標準文字サイズ×TH2」より広いときには、対応した文字矩形ペアの間にスペースが存在すると判断され、対応する認識

結果文字コードの間にスペースコードが挿入される(ステップ510)。同様の処理が対象範囲内の全ての文字矩形間隔比について実行される。

【0031】図7の(a)に示した例では、図7の(d)に示すように、プロポーショナルピッチと判定された対象範囲において「Preter」の末尾の「r」と次の「5」との間の文字矩形間隔比だけがTH2を超えるので、その位置にスペースコードが挿入される。

【0032】定ピッチと判定された対象範囲に対するスペース検出処理(ステップ303)の処理内容について説明する。図6は、そのフローチャートである。

【0033】対象範囲内のすべての文字間の文字矩形ピッチを計算し、結果をワークメモリ212に保存する(ステップ600)。すなわち、文字間を介して隣接する文字矩形の中心を計算し、その中心の間隔を当該文字間の文字矩形ピッチとして求める。この計算には、切り出しメモリ209に保存されている文字矩形の座標データを用いる。ただし、切り出し処理段階で予め隣接した文字矩形の中心間隔を計算して切り出しメモリ209に保存しておき、その値を読み出すようにしてもよい。

【0034】このような文字矩形ピッチの計算が終わると、対象範囲の先頭より、注目する一つの文字間を選び(ステップ605)、注目文字間の前の文字間の文字矩形ピッチと、注目文字間の後の文字間の文字矩形ピッチとを比較し、小さい方の文字矩形ピッチを注目文字間に対する基準ピッチとする(ステップ610)。注目文字間が対象範囲の最初の文字間であるときには、注目文字間の前には文字間が存在しないので、注目文字間の後の文字間の文字矩形ピッチを基準ピッチとする。同様に、注目文字間が対象範囲の最後の文字間であるときには、その前の文字間の文字矩形ピッチを基準ピッチとする。

【0035】次に、注目文字間の文字矩形ピッチを基準ピッチで割った値を文字矩形ピッチ比(=文字矩形ピッチ/基準ピッチ)として計算する(ステップ615)。そして、文字矩形ピッチ比が所定の閾値TH3より大きい、すなわち注目文字間の文字矩形ピッチが「基準ピッチ*TH3」より広いか判定する(ステップ620)。閾値TH3は例えば1.8である(この値に限らないが、一般的文書では、このあたりの値が、全角文字列中の全角スペースも半角文字列中の半角スペースも検出可能であり、適当であることが実験により確認された。ただし、処理する文書に応じて調整するのが好ましい)。

【0036】文字矩形ピッチ比がTH3以下であれば、注目文字間にスペースが存在しないと判断される。文字矩形ピッチ比がTH3より大きいときには、注目文字間にスペースが存在すると判断されるので、結果メモリ211内の対応した文字コードの間にスペースコードが挿入される(ステップ625)。同様の処理が対象範囲の最後の文字間まで繰り返される。

【0037】図7の(a)に示した例では、図7の(d)に示すように、定ピッチと判定された対象範囲において「Imagio」の最後の「o」と次の「7」の文字間だけ文字矩形ピッチ比がTH3を超えるので、その文字間にスペースがあると判断されスペースコードが挿入される。

【0038】

【発明の効果】請求項1記載の発明によれば、文字認識結果より英数文字列を対象範囲として抽出し、対象範囲毎にピッチ書式を判別し、各対象範囲に対して、そのピッチ書式別のスペース検出処理を行うため、行中に英数文字列と日本文字列とが混在し、しかも定ピッチの英数文字列とプロポーショナルピッチの英数文字列が混在するような日英混在文書に対しても、精度のよいスペース検出が可能になる。

【0039】請求項2記載の発明によれば、対象範囲として抽出した英数文字列毎に標準文字サイズを決定し、英数文字列における文字矩形間隔の標準文字サイズに対する相対的な広狭の割合に着目してピッチ書式を判定するため、英数文字列が日本文字列と混在し、しかも全角文字と半角文字の英数文字列が混在するような場合にも、英数文字列のピッチ書式が定ピッチであるかプロポーショナルピッチであるかを精度よく判別でき、したがって日英混在文書中の英数文字列に対し、そのピッチ書式にあったスペース検出処理を的確に適用できるため、定ピッチとプロポーショナルピッチの英数文字列が混在した文書中のスペースを高精度に検出することができる。

【0040】定ピッチ英数文字列の文字矩形ピッチは文字によって大きく変動するため、この変動が反映されない一定の基準ピッチと文字矩形ピッチとの相対的な広狭を調べても、スペースを精度よく検出できない。これに対し、請求項3記載の発明によれば、注目した文字間の前と後の文字間の文字矩形ピッチの中の小さい方を、注目した文字間に対する基準ピッチとして用いるため、文字矩形ピッチの変動が大きい定ピッチ英数文字列中のスペースを精度よく検出することができ、したがって定ピッチの英数文字列が混在した文書に対するスペース検出精度を上げることができる。

【0041】プロポーショナルピッチ英数文字列は文字矩形間隔が狭いため、ピッチ書式や文字種を考慮しない標準文字サイズと文字矩形間隔とを比較する方法では、スペース検出が難しかった。請求項4記載の発明によれば、個々のプロポーショナルピッチ英数文字列毎に標準文字サイズを決定し、標準文字サイズに対する文字矩形間隔の比を閾値処理するため、プロポーショナルピッチの英数文字列中のスペースを精度よく検出でき、したがってプロポーショナルピッチの英数文字列が混在した文書のスペース検出精度を上げることができる。

【0042】請求項5記載の発明によれば、英数文字列

毎に標準文字サイズを決定し、英数文字列における文字矩形間隔の標準文字サイズに対する相対的な広狭の割合に着目してピッチ書式を判定するため、英数文字列が日本文字列と混在し、しかも全角文字と半角文字の英数文字列が混在するような場合にも、英数文字列が定ピッチであるかプロポーショナルピッチであるかを精度よく判別できる。

【0043】定ピッチ英数文字列の文字矩形ピッチは文字によって大きく変動するため、この変動が反映されない一定の基準ピッチと文字矩形ピッチとの相対的な広狭を調べても、スペースを精度よく検出できないが、請求項6記載の発明によれば、注目した文字間の前と後の文字間の文字矩形ピッチの中の小さい方を、注目した文字間に対する基準ピッチとして用いるため、文字矩形ピッチの変動が大きい定ピッチ英数文字列中のスペースを精度よく検出することができる。

【0044】プロポーショナルピッチ英数文字列は文字矩形間隔が狭いため、ピッチ書式や文字種を考慮しない標準文字サイズと文字矩形間隔とを比較する方法では、スペース検出が難しかったが、請求項7記載の発明によれば、個々のプロポーショナルピッチ英数文字列毎に標準文字サイズを決定し、標準文字サイズに対する文字矩形間隔の比を閾値処理するため、プロポーショナルピッチの英数文字列中のスペースを精度よく検出できる。

【図面の簡単な説明】

【図1】本発明の一実施例の全体処理フローを示すフローチャートである。

【図2】本発明の一実施例のための装置構成の例を示すブロック図である。

【図3】図1中のスペース検出処理の概略を示すフローチャートである。

【図4】図3中のピッチ書式判定処理の内容を示すフローチャートである。

【図5】図3中のプロポーショナルピッチ文字列スペース検出処理の内容を示すフローチャートである。

【図6】図3中の定ピッチ文字列スペース検出処理の内容を示すフローチャートである。

【図7】スペース検出の具体例を説明するための図である。

【図8】ピッチ書式判定のための標準文字サイズの決定方法の説明図である。

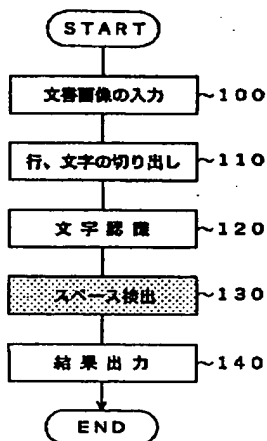
【図9】日英混在文書のスペース検出における課題を明らかにするための図である。

【符号の説明】

- 200 画像入力部
- 202 行・文字切り出し部
- 203 文字認識部
- 204 対象範囲抽出部
- 205 ピッチ書式判定部
- 206 定ピッチ文字列スペース検出部
- 207 プロポーショナルピッチ文字列スペース検出部
- 208 画像メモリ
- 209 切り出しメモリ
- 210 文字辞書メモリ
- 211 結果メモリ
- 212 ワークメモリ
- 213 バス

【図1】

全体処理フロー



【図8】

標準文字サイズ決定の説明図

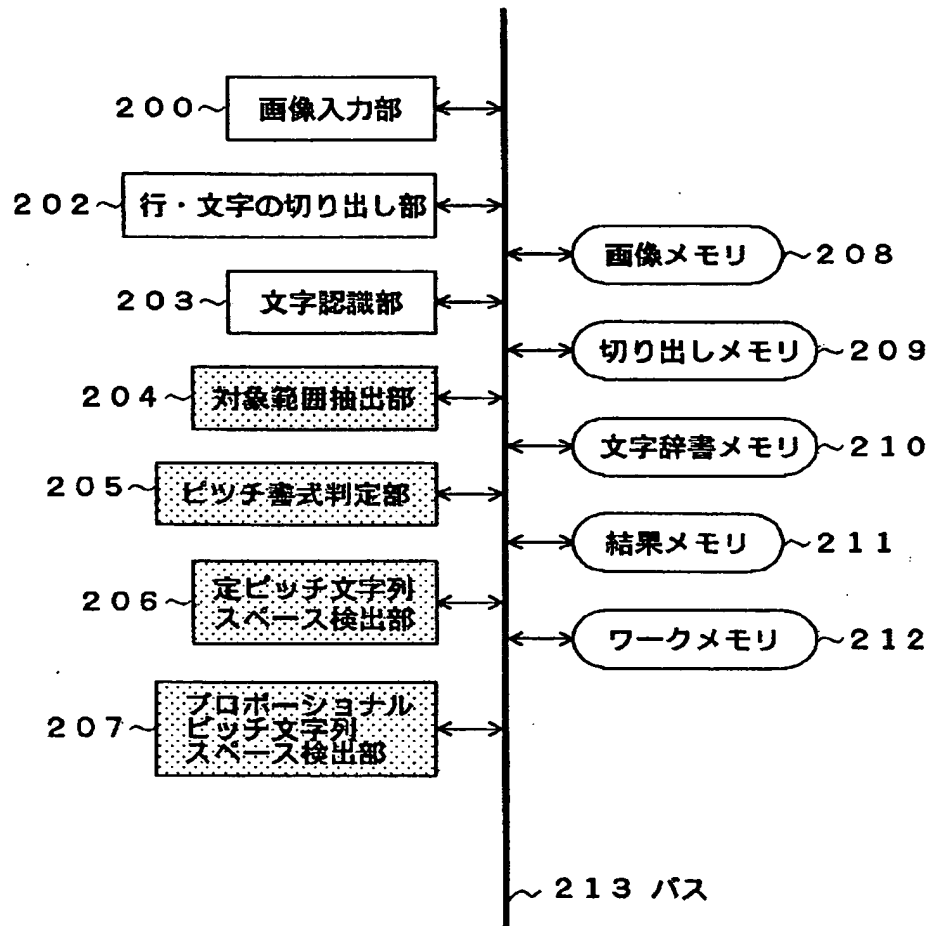
対象文字列例:

漢字 ひら Abeg Abcg カタ列
対象範囲

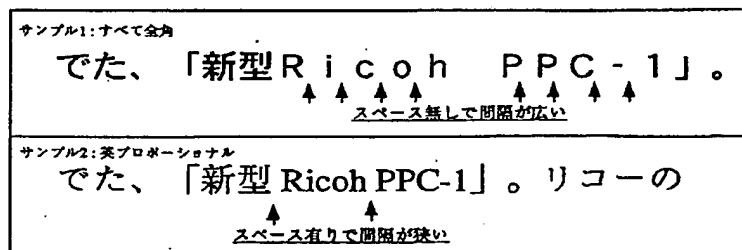
漢 び 力 A B

文字高さの比較: $\left[\frac{\text{漢の文字高さ}}{\text{ひの文字高さ}} \right] \approx 1.25$ 標準文字サイズ

【図2】

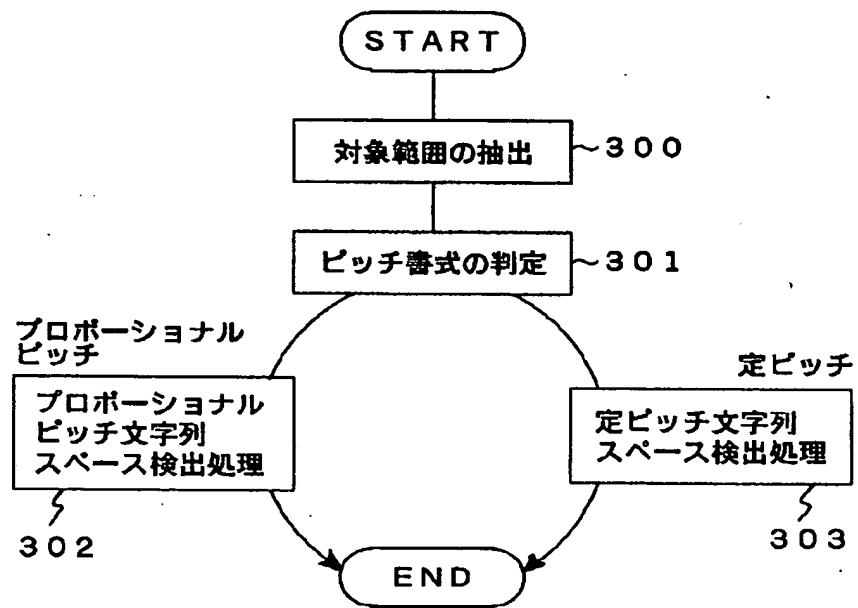


【図9】



【図3】

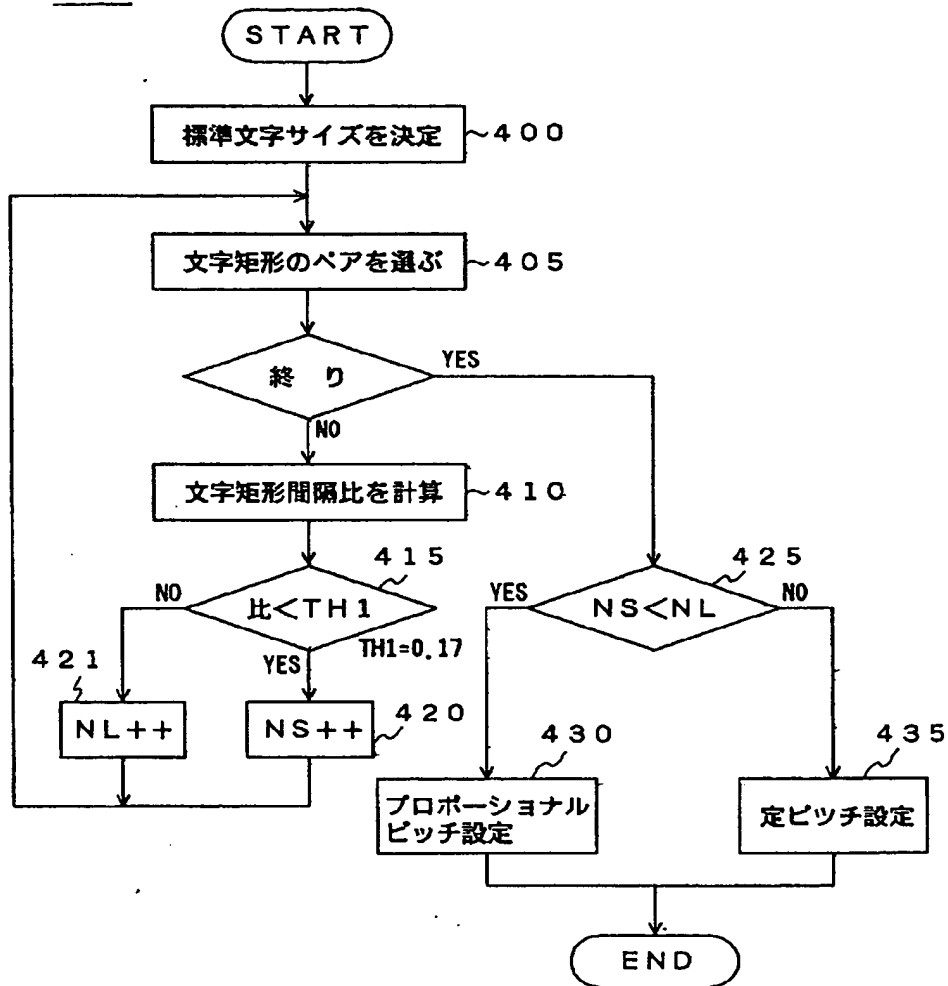
スペース検出処理フロー



【図4】

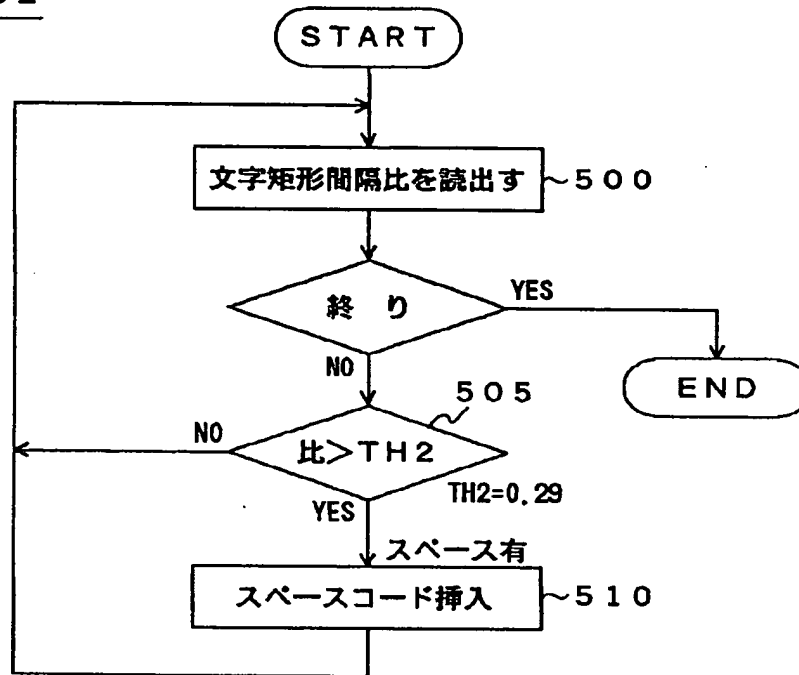
ピッチ書式判定

301



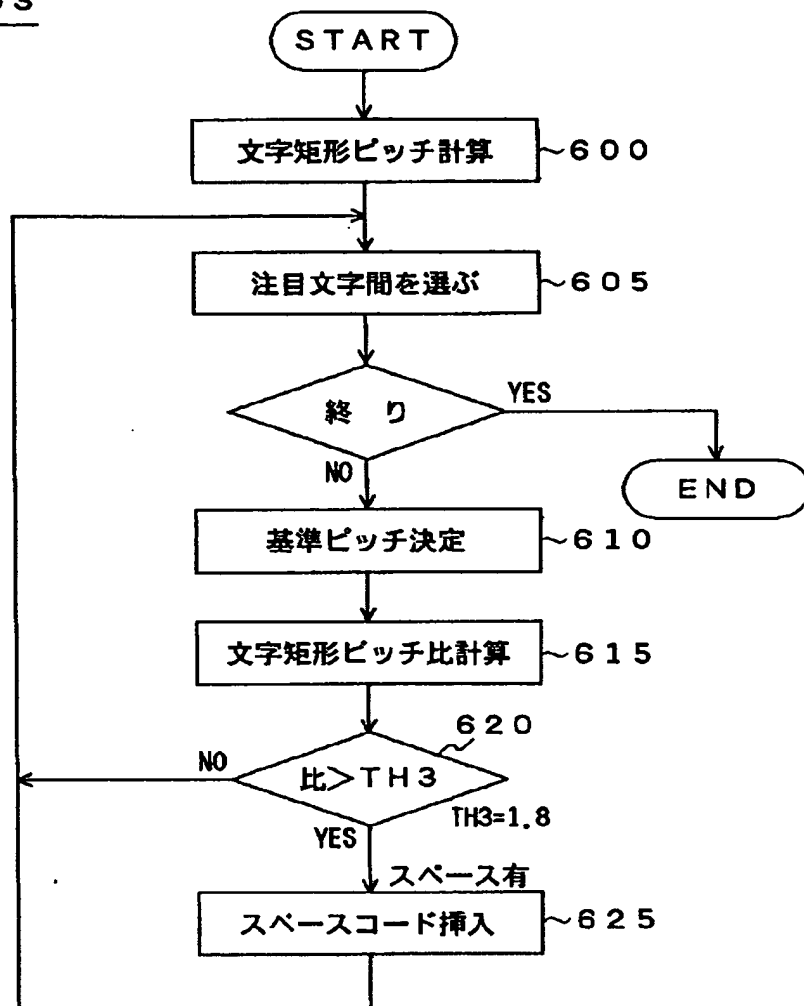
【図5】

プロポーショナルピッチのスペース検出

302

【図6】

定ピッチのスペース検出

303

【図7】

スペース検出の処理例

(a) フォントの混在する日英混在文字列サンプル

新型コピー、Preter 55とImagio 77発売]

(b) 対象範囲の抽出

・後処理後結果5文字以上
連続する英数文字列

---(Times-Roman)---

------(Courier)-----

行最大
文字矩形
高さ

8文字

8文字

(c) ピッチ書式判定

英数最大
文字矩形
高さ
[Preter 55]
文字矩形間隔英数最大
文字矩形
高さ
[Imagio 77]
文字矩形間隔標準文字サイズ = $\text{Max}(\text{行最大高さ}, \text{英数最大高さ} \times 1.25)$
文字矩形間隔 < 標準文字サイズ * 0.17

文字矩形間隔狭い、が多数

プロポーショナルピッチ

標準文字サイズ = $\text{Max}(\text{行最大高さ}, \text{英数最大高さ} \times 1.25)$
文字矩形間隔 \geq 標準文字サイズ * 0.17

文字矩形間隔狭い、が多数

定ピッチ

(d) スペース判定

[Preter 55]
文字矩形間隔スペース判定しきい値
標準文字サイズ * 0.29 より大きい

文字間スペース: 無 無 無 無 無 有 無

[Imagio 77]
文字矩形間
ピッチ自分の両側のピッチの小さい方の
1.8倍よりも自分が大きい

文字間スペース: 無 無 無 無 無 有 無

